

The World's Largest Open Access Agricultural & Applied Economics Digital Library

This document is discoverable and free to researchers across the globe due to the work of AgEcon Search.

Help ensure our sustainability.

Give to AgEcon Search

AgEcon Search http://ageconsearch.umn.edu aesearch@umn.edu

Papers downloaded from **AgEcon Search** may be used for non-commercial purposes and personal study only. No other use, including posting to another Internet site, is permitted without permission from the copyright owner (not AgEcon Search), or as allowed under the provisions of Fair Use, U.S. Copyright Act, Title 17 U.S.C.



STOCHASTIC GLOBAL OPTIMIZATION METHODS

PART I: CLUSTERING METHODS

A.H.G. Rinnooy Kan* ** G.T. Timmer** ***

ABSTRACT

In this stochastic approach to global optimization, clustering techniques are applied to identify local minima of a real valued objective function that are potentially global. Three different methods of this type are described; their accuracy and efficiency are analyzed in detail.

* Department of Industrial Engineering and Operations Research/Graduate School of Business Administration, University of California, Berkeley

** Econometric Institute, Erasmus University Rotterdam

*** ORTEC Consultants, Rotterdam

1. INTRODUCTION

Let $f : \mathbb{R}^n \to \mathbb{R}$ be a continuous real valued objective function. Most nonlinear programming methods that have been developed aim for a <u>local</u> <u>optimum</u> (say, local minimum), i.e. a point x* such that there exists a neighbourhood B of x* with

 $f(x^*) \leq f(x) \quad \forall x \in B$ (1)

In general, however, several local optima may exist and the corresponding function values may differ substantially. The problem of designing algorithms that distinguish between these local optima and locate the best possible one is known as the <u>global optimization</u> problem, and forms the subject of this paper and its companion, Part II.

In the absence of reliable codes for the global optimization problem most problems are not modelled as such. Many problems, however, are of a global nature. This is especially true for many <u>technical design</u> problems [Dixon & Szegö 1978b, Archetti & Frontini 1978]. <u>Economic applications</u>, where multimodal cost functions have to be minimized, have also been reported [Archetti & Frontini 1978]. Another global optimization problem often encountered in econometrics is that of locating the global maximum of a likelihood function. Thus, there is no need to dwell on the practical usefulness of quick and reliable methods to solve the global optimization problem.

The global optimization problem is to find the <u>global optimum</u> (say global minimum) x_* of a real valued objective function $f : \mathbb{R}^n \to \mathbb{R}$, i.e. to find a point $x_* \in \mathbb{R}^n$ such that

 $f(x_*) \leq f(x) \quad \forall x \in \mathbb{R}^n.$ (2)

Unless stated otherwise, we will assume f to be twice continuously differentiable. For obvious computational reasons, one usually assumes that a set $S \subset \mathbb{R}^n$, which is convex, compact and contains the global minimum as an interior point, is specified in advance. None the less, the problem to find

$$y_* = \min_{x \in S} f(x)$$

remains essentially one of unconstrained optimization.

Any method for global optimization has to account for the fact that a numerical procedure can never produce more than approximate answers. Thus, the global optimization problem might be considered solved if, for some $\varepsilon > 0$, an element of one of the following sets has been identified [Dixon 1978]

$$A_{\mathbf{x}}(\varepsilon) = \{\mathbf{x} \in S \mid \|\mathbf{x} - \mathbf{x}_{\mathbf{x}}\| \leq \varepsilon\},$$
(4)

$$A_{f}(\varepsilon) = \{x \in S | |f(x) - f(x_{*})| \leq \varepsilon \}.$$
(5)

A disadvantage of the first mentioned possibility is that small perturbations in the problem data may have major effects on the location of x_{\star} [Archetti & Betro 1978a]. A third possibility [Betro 1981] is obtained by defining

$$\phi(\mathbf{y}) = \frac{m(\{z \in S \mid f(z) \leq \mathbf{y}\})}{m(S)}, \qquad (6)$$

where m(.) is the Lebesque measure and taking

$$A_{\phi}(\varepsilon) = \{x \in S | \phi(f(x)) \leq \varepsilon\}.$$
(7)

We note, however, that this set may contain points whose function values differ considerably from y_* .

A second problem, which is caused by the finite accuracy of numerical procedures, is that we cannot distinguish between two local minima which are very close to one another. If we define a stationary point of f as a point where the gradient g : $\mathbb{R}^n \to \mathbb{R}^n$ of f is equal to 0, then each (local) minimum is known to be a stationary point. We will assume that a positive constant ε can be specified, such that the distance between any two stationary points exceeds ε . Obviously, this implies that there can only be a finite number of stationary points in S.

-

(3)

Only few solution methods for global optimization have been developed so far; we refer to [Dixon & Szegö 1978a, 1978b] and to [Rinnooy Kan & Timmer 1984, Boender et al. 1985] for surveys. We shall be concerned with methods that incorporate <u>stochastic</u> elements. In most stochastic methods, two phases can be usefully distinguished. In the <u>global phase</u>, the function is evaluated in a number of randomly sampled points. In the <u>local phase</u>, the sample points are manipulated, e.g. by means of local searches, to yield a candidate global minimum.

Generally in turning to stochastic methods, we do sacrifice the possibility of an <u>absolute guarantee</u> of success. However, under mild conditions on the sampling distribution and on f, the probability that an element of $A_x(\varepsilon)$, $A_f(\varepsilon)$ or $A_\phi(\varepsilon)$ is sampled approaches 1 as the sample size increases [Solis & Wets 1981]. If the sample points are drawn from a <u>uniform</u> distribution over S and if f is continuous, then an even stronger result will turn out to hold: the sample point with lowest function value converges to the global minimum value with <u>probability 1</u> (or almost surely). Thus, the global phase can yield an asymptotic guarantee with probability 1, and is therefore essential for the <u>reliability</u> of the method. However, a method that only contains a global phase will be found lacking in <u>efficiency</u>. To increase the latter while maintaining the former is one of the challenges in global optimization.

As in the case of deterministic methods, one of the questions in applying a stochastic method is when to stop. Preferably, a method of this nature should terminate with some probabilistic information on the quality of the proposed solution. Several approaches based on different assumptions about the properties of possible objective functions f and using different stochastic techniques have been proposed to design a proper <u>stopping rule</u>.

In Section 2, we review some stochastic methods and find that the most promising methods appear to be variants of the so-called <u>Multistart</u> technique where points are sampled iteratively from a <u>uniform distribution</u> over S (global phase), after which local minima will be found by applying a <u>local search procedure</u> to these points (local phase). A theoretical framework which enables the stochastic analysis of this method is developed in [Boender 1984] (see also [Boender & Rinnooy Kan 1983, 1985]). It turns out to be possible to develop <u>Bayesian estimates</u> of the number of local minima not yet identified and of the probability that the next local search will locate a new local minimum. By specifying the costs and the potential

benefits of further experiments and weighing these against each other probabilistically, an <u>optimal Bayesian stopping rule</u> can be determined.

Multistart is still lacking in efficiency because the same local minimum may be located several times. If we define the <u>region of attraction</u> of a local minimum x* to be the set of points in S starting from which a given local search procedure converges to x*, then ideally, this local search procedure should be started exactly once in every region of attraction. Several new algorithms designed to satisfy this criterion are described in Section 3. The methods discussed in this section temporarily eliminate a prespecified fraction of the sample points whose function values are relatively high. The resulting <u>reduced sample</u> consists of groups of mutually relatively close points that correspond to the regions with relatively small function values. Within each group the points are still distributed according to the original uniform distribution. Thus, these groups can be identified by <u>clustering techniques</u> based upon tests on the uniform distribution. Only one local search procedure will be started in each group [Boender et al. 1980, 1982].

Unfortunately, the resulting groups do not necessarily correspond to the regions of attraction of f. It is possible that a certain group of points corresponds to a region with relatively small function values which contains several minima. Therefore, the methods which are based on the reduced sample may fail to find a local minimum although a point is sampled in its region of attraction. Methods that do not suffer from this deficiency will be dealt with in Part II of this paper [Rinnooy Kan & Timmer 1985]. There we also discuss the computer implementation of the various global optimization methods and its theoretical properties, and we discuss the results of some computational experiments.

MULTISTART

The simplest stochastic method for global optimization consists only of a global phase. Known confusingly as Pure Random Search [Brooks 1958, Anderssen 1972], the method involves no more than a single step.

Pure Random Search

Step 1. Evaluate f in N points, drawn from a uniform distribution over S. The smallest function value found is the candidate solution for y*.

In spite of its evident simplicity, Pure Random Search offers an <u>asymptotic guarantee</u> in a probabilistic sense. The proof is based on the simple observation that the probability that a uniform sample of size N contains at least one point in a subset $A \subset S$ is equal to [Brooks 1958]

$$1 - (1 - \frac{m(A)}{m(S)})^{N}$$
. (8)

Thus, any assumption on f guaranteeing that $m(A_f(\varepsilon))$, $m(A_x(\varepsilon))$ or $m(A_{\phi}(\varepsilon))$ is strictly positive will imply that Pure Random Search locates an element in the corresponding set with a probability approaching to 1 as N increases. In fact, if we let $\underline{y}_N^{(1)}$ be the smallest function value found in a sample of size N, then we can prove the following result.

THEOREM 1. (cf. [Devroye 1978, Rubinstein 1981]) If f is continuous, then $\underline{y}_{N}^{(1)}$ converges to the global minimum value y_{\star} with <u>probability 1</u> (or <u>almost surely</u>) with increasing N, i.e.

 $\Pr\left[\lim_{N \to \infty} \underline{y}_{N}^{(1)} = y_{\star}\right] = 1.$ (9)

<u>PROOF</u>. Because f is continuous in a global minimum x_* , we know that for all $\varepsilon > 0$ there exists a $\delta > 0$, such that the probability that $\underline{y}_N^{(1)} - y_* > \varepsilon$ is less than the probability that no element in a sample of size N is within distance δ from that global minimum. Since x_* is assumed to be in the interior of S, we can choose δ small enough so as to ensure that the set of points which are within distance δ of x_* is completely contained in S. Hence

$$\Pr\left[\left|\underline{y}_{N}^{(1)}-y_{\star}\right| > \varepsilon\right] \leq (1-\delta)^{N}, \qquad (10)$$

and

$$\sum_{N=1}^{\infty} \Pr\left[\left|\underline{y}_{N}^{(1)} - y_{\star}\right| > \varepsilon\right] \leq \sum_{N=1}^{\infty} (1 - \delta)^{N} = \frac{1 - \delta}{\delta}.$$
 (11)

Thus, the left-hand side of (11) converges for all ε , so that, using the Borel-Cantelli Lemma [Chung 1974], (9) follows immediately.

A similar guarantee will hold for all methods that follow.

The reader may well wonder to what extent an embarassingly simple method such as Pure Random Search has any advantage over an equally simplictic approach such as Grid Search, in which the function is evaluated in each point of a regular grid over S. The relative merits of these naive stochastic and deterministic strategies have been extensively analyzed [Sukharev 1971, Ivanov 1972, Anderssen & Bloomfield 1975, Archetti & Betro 1978b]. The net result of these analyses is that the points of the random sample cover S more efficiently (according to several probabilistic criteria) than the grid points do, at least if the dimension of the problem is not too low. In the studies mentioned above the methods are evaluated according to the distance between the global minimum and the sample or grid point closest to it. The advantage of Pure Random Search becomes more evident through an argument in [Sobol 1982]. Here, it is observed that for many functions some of the variables (in the n-dimensional space S) hardly affect the function value; in which case the distribution of the sample or grid points in the subspace defined by the remaining variables is of primary interest. However, it is not known in advance which of the variables are important and which are not. If the (uniform) sample points are projected into an arbitrary subspace, they still follow a uniform distribution over this subspace. However, if the grid points are projected into an arbitrary subspace, they may very well form groups of mutually close points, that cover the subspace in an unsatisfactory manner.

None the less, Pure Random Search can hardly be taken seriously as a computational proposal. Several extensions of this method have been proposed that also start from a uniform sample over S (hence, Theorem 1 can be applied), but that at the same time involve local searches from some or all points in the sample. The simplest way to make use of a <u>local search</u> procedure P occurs in a folklore method known as Multistart.

Multistart

- Step 1. Draw a point from a uniform distribution over S.
- Step 2. Apply P to the new sample point.
- Step 3. A termination criterion indicates whether to stop or to return to Step 1. The local minimum with smallest function value found is the candidate value for y_* .

Although this method is obviously more attractive then Pure Random Search several inefficiencies still remain. However, let us first consider the question of a proper stopping rule for this method. Our treatment will be brief, since the details of our approach are reported elsewhere; it was initiated in [Zielinski 1981] and extended in [Boender 1984]. It is based on a Bayesian estimate of the number of local minima W and the relative size of each region of attraction $\Theta_{\ell} = m(R_{x*})/m(S)$, $\ell=1,...,W$, where R_{x*} is the region of attraction of the local minimum x*, i.e., the set of points in S starting from which P will converge to x*. If the values of these parameters would be given, then the outcome of an application of Multistart is easy to analyze. We can view the procedure as a series of experiments in which a sample from a multinomial distribution is taken. Each <u>cell</u> of the distribution corresponds to a minimum x*; the <u>cell</u> probability is equal to the probability that a uniformly sampled point will be allocated to $R_{_{\bf X}\bigstar}$ by P, i.e. equal to the corresponding $\Theta_{{\it l}}.$ Thus, the probability that the l-th local minimum is found b_l times ($l = 1, \dots, W$) in N trials is

$$\frac{N!}{\underset{\substack{\mathbb{N} = 1}}{\mathbb{W}}} \cdot \underset{\substack{\ell=1}}{\overset{\mathbb{W}}{\mathbb{H}}} \Theta_{\ell}^{\mathbf{b}_{\ell}}.$$
 (12)

It is impossible, however, to distinguish between outcomes that are identical up to a relabeling of the minima. Thus, we have to restrict ourselves to <u>distinguishable aggregates</u> of the random events that appear in (12). To calculate the probability that w different local minima are found during N local searches, and the i-th minimum is found a_i times $(a_i > 0, \sum_{i=1}^{W} a_i = N)$, let c_j be the number of a_i 's equal to j and let $S_W(w)$ denote the set of all permutations of w different elements of $\{1, 2, \dots, W\}$. The required probability is then given by [Boender 1984]

$$\frac{1}{N} \cdot \frac{N!}{w} \cdot \sum_{\substack{w \\ i=1}}^{w} (\pi_1, \dots, \pi_w) \in S_W(w) \xrightarrow{w}_{i=1}^{u} \Theta_{\pi_i}^{i}.$$
(13)

Formula (13) can be used in a <u>Bayesian</u> approach in which the unknowns W, $\Theta_1, \ldots, \Theta_W$ are assumed to be themselves random variables for which a <u>prior distribution</u> can be specified. Given the outcome of an application of Multistart, Bayes's rule is used to compute the <u>posterior</u> <u>distribution</u>, which incorporates both the prior beliefs and the sample information.

After lengthy calculations, surprisingly simple expressions emerge for the posterior distribution and posterior expectation of several interesting parameters, some of which are stated in the next theorem.

<u>THEOREM 2.</u> ([Boender 1984]) If w different local minima have been found as the result of N local searches started in uniformly distributed points, if we assume a priori for the number of local minima <u>W</u> that each integer of $[1,\infty)$ is equally probable, and if we assume that given <u>W</u> = W the relative sizes of the regions of attraction $\underline{\Theta}_1, \dots, \underline{\Theta}_W$ follow a uniform distribution on the (W-1)-dimensional unit simplex, then

i) the posterior probability that there are K local minima is equal to

$$\frac{(K-1)!K!(N-1)!(N-2)!}{(N+K-1)!(K-w)!w!(w-1)!(N-w-2)!}$$
(14)

ii) the posterior expectation of the number of local minima is

$$\frac{w(N-1)}{N-w-2}$$
(15)

iii) the posterior expected size of the non-observed regions of attraction is

$$\frac{w(w+1)}{N(N-1)}$$
(16)

This theoretical framework is quite an attractive one, the more so since it can be easily extended to yield <u>optimal Bayesian stopping rules</u>. As in the previous section, such rules incorporate assumption about the costs and potential benefits of further experiments and weigh these against each other probabilistically to calculate the optimal stopping point. Several loss structures and corresponding stopping rules are described in [Boender & Rinnooy Kan 1985].

None the less, in spite of the scope that Multistart offers for analysis, the procedure is still lacking in efficiency. The main reason for this is that it will inevitably cause each local minimum to be found several times. To avoid all these time consuming local searches, P should ideally be invoked no more than once in every region of attraction.

A first attempt to modify Multistart in this way can be found in [Hartman 1973]. In this method, a local search is started only when a point is drawn whose function value is less than the smallest local minimum value found so far. It should be obvious that under this rule the global minimum may not be found even if a point is sampled in R_{x*} . A far more successful adaptation of Multistart is provided by the <u>clustering methods</u>, which form the main subject of this paper.

The basic idea behind the clustering methods is to start from a uniform sample from S, to create groups of mutually close points that correspond to relevant regions of attraction, and to start P once in every such region. Two ways to create such groups from the initial sample have been proposed. The first, called reduction [Becker & Lago 1970] removes a certain fraction of the sample points with the highest function values. The second, called <u>concentration</u> [Törn 1978], transforms the sample by allowing one or at most a few <u>steepest descent</u> steps from every point.

A disadvantage of concentration is that it transforms the sample in an

unpredictable way. Therefore, the methods based on concentration may fail in two different ways. Firstly, the resulting groups of points, or <u>clusters</u>, may contain several regions of attraction, so that the global minimum can be missed. Secondly, one region of attraction may be divided over several clusters, in which case the corresponding minimum will be located more than once. A further disadvantage of concentration is that nothing can be said about the distribution of the resulting points, which makes it more difficult to identify the clusters. Better results will be seen to be possible for methods which are based on reduction. These methods will be dealt with in the next section.

3. CLUSTERING METHODS

In this section we aim for solution methods for this global optimization problem that satisfy the following - partially conflicting demands. On the one hand the method must be <u>asymptotically correct</u>, i.e. if the method would be continued sufficiently long, then the smallest function value found during this process must converge to the global minimum value. (If the method is stochastic, convergence in a probabilistic sense will be required.) On the other hand, the method must be <u>efficient</u>: if the method is stopped after a reasonable amount of time, it should have produced results which compare favourably to the results obtainable by other methods in the same time period.

For reasons given in the previous section, the methods we will consider are variants of Multistart. The methods are iterative, and fit in the following framework.

Global framework

- Step 1. (Global phase) N points are drawn from a uniform distribution over S. The function is evaluated in these points, and the points are added to the (initially empty) sample.
- <u>Step 2.</u> (Local phase) A procedure selects a (possibly empty) subset of the enlarged sample, and a local search procedure P is applied to each of the elements of this subset. The stationary points, which are found during these local searches and which were not detected previously, are added to an (initially empty) set X^{*}.
- <u>Step 3</u>. A stopping rule decides whether to return to Step 1 or to stop. If the method is stopped, then the element of X^{*} with smallest function value is the candidate solution.

In this section, we assume that the local search procedure P in Step 2 is strictly descent: starting from any point $x \in S$, it generates a sequence of points x_k , with

$$x_{k+1} = x_k + \alpha_k p_k \quad (\|p_k\| = 1, \alpha_k > 0), \quad (17)$$

which converges to a stationary point x, such that moreover

$$f(x_k + \beta p_k) \leq f(x_k + \alpha p_k)$$
(18)

for all k and all α , β satisfying $0 \leq \alpha < \beta \leq \alpha_k$. Thus, there exists a path from x to \overline{x} along which the function is nonincreasing. We shall also assume that this path is completely contained in S. As a result, we can now derive some important properties of the regions of attraction of P.

Let

$$L(y) = \{x \in S \mid f(x) \leq y\}$$
(19)

be the (y-) <u>level set</u> of f, and let $L_x(y)$ (with $y \ge f(x)$) denote the (connected) <u>component</u> of L(y) containing x. It is not difficult to see that both these sets contain all their accumulation points and are hence closed.

<u>Theorem 3.</u> If, for any $x \in S$ and $y \geq f(x)$, a procedure which is strictly descent is started from a point in $L_x(y)$, then the sequence x_k generated by this procedure will converge to a stationary point \overline{x} in $L_x(y)$.

<u>PROOF</u>. Since all points of the sequence x_k are located in S and since S is convex, the interval $[x_k, x_{k+1}]$ is completely contained in S, for every k. Suppose that there exists a k such that $x_k \in L_x(y)$ and $x_{k+1} \notin L_x(y)$. It follows that there exists an element \overline{x} on the line segment $[x_k, x_{k+1}]$ with $f(\overline{x}) > y$. Because, if there would not be such an element, then this line segment would be contained in L(y) and since the line segment is clearly path-connected (and therefore connected), it would follow from the definition of a component, that x_k and x_{k+1} belong to the same component of L(y). However, if $f(\overline{x}) > y$, then $f(\overline{x}) > f(x_k)$ which contradicts (18). We conclude that $x_k \in L_x(y)$ for every k. The sequence x_k must converge to a stationary point \overline{x} . Since $x_k \in L_x(y)$ for every k and $L_x(y)$ is closed, the result is immediate.

For some local minimum x*, let $\overline{y} \in \mathbb{R}$ be the largest y for which the interior of $L_{x*}(y)$ contains x* as its only stationary point (because there

13

are only a finite number of stationary points, this \overline{y} really exists). We define the <u>basin</u> B_{x*} of x* as the interior of $L_{x*}(\overline{y})$.

THEOREM 4. ([Dixon et al. 1975]) If a procedure that is strictly descent is started from any point x in B_{x*} , then the sequence x_k generated by this procedure converges to x*, i.e. $\subset R_{x*}$ for strictly descent procedures.

<u>PROOF</u>. For any $x \in B_{x^*}$, $L_{x^*}(f(x))$ contains x^* as its only stationary point. Through Theorem 3 the result is now immediate.

We note that, although condition (18) cannot be verified computationally a slight variation of it is more tractable. Let us define a procedure satisfying (17) to be $\underline{\varepsilon}$ -descent if the sequence converges to a stationary point; if $f(x_{k+1}) \leq f(x_k)$ for all k and moreover

$$f(x_k + i\varepsilon p_k) \leq f(x_k + (i-1) \varepsilon p_k) \quad (i = 1, 2, \dots, [\frac{\alpha_k}{\varepsilon}]). \quad (20)$$

Results similar to Theorem 3 and 4 hold for ε -descent procedures.

<u>THEOREM 5.</u> If, for some $x \in S$ and $y \geq f(x)$, there is no point x_1 with $x_1 \in L(y)$, $x_1 \notin L_x(y)$ which is within ε -distance of an element $L_x(y)$ (i.e. $L_x(y)$ is not too close to another component of L(y)), then an ε -descent procedure started from a point in $L_x(y)$ will converge to a stationary point in $L_x(y)$.

<u>PROOF</u>. The sequence x_k generated by the procedure converges to a stationary point \bar{x} . Suppose that there exists a k such that $x_k \in L_x(y)$ and $x_{k+1} \notin L_x(y)$. Because the procedure cannot leave S, $x_{k+1} \in L(y)$. Since the distance between $L_x(y)$ and any other component of L(y) is known to exceed ε , the line segment $[x_k, x_{k+1}]$ must contain an interval of length ε on which the function values exceed y. At least one of the points $x_k + i\varepsilon p_k$ $(i = 1, 2, \dots, [\frac{\alpha}{\varepsilon}])$ will be located on this interval, so that condition (20) is not satisfied. Hence, $x_k \in L_x(y)$ for every k, and since $L_x(y)$ is closed, $\bar{x} \in L_x(y)$.

For any local minimum x*, we define

$$\begin{split} B_{x^{\star}}^{\varepsilon} &= \left\{ x \in B_{x^{\star}} \right| \exists x_{1} \text{ with } x_{1} \notin B_{x^{\star}} \text{ and } \|x - x_{1}\| \leq \varepsilon \right\}, \\ y_{\varepsilon} &= \inf_{x \in B_{x^{\star}}^{\varepsilon}} f(x), \\ B_{x^{\star}}(\varepsilon) &= \left\{ x \in B_{x^{\star}} \right| f(x) < y_{\varepsilon} \right\}. \end{split}$$

i.e. $B_{X^*}(\varepsilon)$ is the subset of B_{X^*} which consist of points that have a smaller function value than any point in B_{X^*} that is within ε -distance of the boundary of B_{X^*} .

THEOREM 6. If a procedure that is ε -descent is started from any point x in $B_{x*}(\varepsilon)$, then the sequence x_k generated by this procedure converges to x^* .

<u>PROOF.</u> If $B_{x*}(\varepsilon)$ is empty, the theorem is clearly true. If $B_{x*}(\varepsilon)$ is not empty, then it containts the local minimum x*. For any k, suppose that $x_k \in B_{x*}(\varepsilon)$ and $x_{k+1} \notin B_{x*}(\varepsilon)$. Then there are two possibilities: either $x_{k+1} \in B_{x*}$ (and $x_{k+1} \notin B_{x*}(\varepsilon)$), or $x_{k+1} \notin B_{x*}$. The first possibility cannot occur because it would imply that $f(x_{k+1}) > f(x_k)$. The second possibility cannot occur, since the line segment $[x_k, x_{k+1}]$ would then contain an interval of length ε on which the function values exceed $f(x_k)$, in which case (20) is not satisfied for all i. We conclude that $x_k \in B_{x*}(\varepsilon)$ for all k. Since x* is the only stationary point in the closure of $B_{x*}(\varepsilon)$ (x* is the only stationary point B_{x*} and the closure of $B_{x*}(\varepsilon)$ is a subset of B_{x*}), the procedure converges to x*.

As a final assumption, let us suppose that \bar{x} , the ultimate point of convergence for P, is actually a local minimum of f. This assumption seems to be an innocent one from an empirical point of view [Wolfe 1969, 1971]; if P should get stuck in a saddlepoint, we could always leave it after a suitable pertubation.

let us now return to the global framework mentioned above and assume that P is strictly descent. If the subset selected in Step 2 equals the set of points which are added to the sample in Step 1, then the global

framework reduces to Multistart. It is not efficient, however, to apply P to every sample point. Obviously, it would be preferable to apply P to a sample point if and only if this point is located in a region of attraction belonging to a minimum that has not yet been found. The set of minima found would then equal the set of minima found by Multistart, but it would probably be obtained at less costs. The purpose of our research has been to develop a method in which P is started exactly once in every region of attraction in which points have been sampled.

We will first examine the case in which P is only applied to sample points with relatively small function value. More precisely, we will consider procedures that start by temporarily removing a prespecified fraction 1- γ of the sample points ($0 < \gamma < 1$), whose function values are relatively high. The remaining points form the <u>reduced sample</u>. If $y_k^{(i)}$ is the i-th smallest function value in a sample of size kN (obtained after k iterations of the global framework), then all elements of the reduced sample are element of $L(y_k^{(\gamma kN)}) = \{x \in S \mid f(x) \leq y_k^{(\gamma kN)}\}$. (Let us note here that, to facilitate the notation, we shall ignore various necessary integer round-ups and round-downs as in the case of γkN ; they do not affect the analysis at all.) Again it is not very efficient to actually apply P to every reduced sample point, i.e. every point in $L(y_k^{(\gamma kN)})$. Instead, we will seek for methods in which P is started exactly once in every region of attraction which contains a reduced sample point.

We now examine the consequences of this approach for the stopping rule in Step 3, which decides whether or not the search for the global minimum has been sufficiently thorough. Recall that the Bayesian stopping rules described in Section 2 only depend on the number of points sampled and the number of minima that are obtained by starting local searches at these points. Therefore, they are not only applicable to Multistart, but to every method which, given a sample, results in the same set of minima as Multistart. In particular, these stopping rules are applicable to the methods in which exactly one local search is started in every region of attraction in which points have been sampled.

For methods in which P is applied exactly once in every region of attraction containing at least one reduced sample point, the situation is more complicated. To analyze this situation, for any γ with $0 < \gamma < 1$, let

 $y_{\gamma} \in \mathbb{R}$ be such that

$$\phi(y_{\gamma}) = \frac{m(\{x \in S \mid f(x) \leq y_{\gamma}\})}{m(S)} = \gamma, \qquad (21)$$

i.e., y_{γ} is the <u> γ -quantile</u> of f. Since ϕ is a monotonically increasing continuous function, there exists a unique value y_{γ} satisfying (21). If we would apply P to every sample point in $L(y_{\gamma})$, then the Bayesian analysis, as described in Section 2, can still be applied. We can simply ignore the sample points whose function value exceeds y_{γ} , and apply the Bayesian analysis to the remaining points. The analysis can then be adapted to this new situation in a trivial way.

However, since we do not know y_{γ} in advance we cannot apply P to the sample points in $L(y_{\gamma})$. Instead, we aim for methods in which P is applied to points in the level set $L(y_k^{(\gamma kN)})$, such that all minima whose regions of attraction contain a reduced sample point are found. Hence, the level above which the sample points are ignored depends on the sample. Therefore, the cell probabilities are no longer constant over time and the Bayesian analysis is formally no longer applicable. However, it is known that $\underline{y}_k^{(\gamma kN)}$ does converge to y_{γ} with probability 1 [Bahadur 1966]. Hence, we may apply the adapted stopping rules as though $\underline{y}_k^{(\gamma kN)}$ does not vary with k.

Unfortunately, we will not succeed entirely in our search for a method in which P is started exactly once in every region of attraction which contains a sample point, respectively a reduced sample point. In particular, we will not be able to exclude the possibility that P is not applied to a (reduced) sample point, although it would have led to a local minimum which has not yet been found. However, we apply the stopping rules as though this possibility does not exist, and will justify our use of these rules by showing that the probability that an error of the above type is made goes to 0 when the sample size increases.

Now, given kN sample points that have been drawn from a uniform distribution over S and given a set of stationary points X^* , we must determine a subset of the sample points to which P will be applied. To do so, we will use the reduced sample to estimate the components of $L(y_k^{(\gamma kN)})$. A local search is then started once in each component that does not contain an element of X^* . The rationale of this approach is that if P, is applied to an element of a component of $L(y_k^{(\gamma kN)})$, then P is known

to converge to a local minimum in that component (see Theorem 3).

How can the components of $L(y_k^{(\gamma kN)})$ be identified? Intuitively, as a result of the removal of a fraction of the points with higher function values, groups of points that are relatively close to each other are created, each of which corresponds to such a component. The natural way to identify these groups (and through them, the components) is to make use of <u>cluster analysis</u>. However, there are several reasons not to use the ordinary clustering methods. The main reason is that we have more information about the problem than just the location and the function value of the reduced sample points. This extra information includes the fact that the reduced sample points are known to be a subset of a uniform sample and the fact that the groups searched for correspond to the components of a level set of a continuously differentiable function. Since this information cannot be translated into measurable characteristics of the reduced sample points, it must be ignored or used in a different way.

The methods which we will describe shortly can be viewed as standard clustering techniques which have been adapted to our specific problem, and all fit in the following framework. The clusters are created one by one, and each cluster is initiated by a <u>seed point</u>. Selected points of the reduced sample are added to the cluster until a <u>termination criterion</u> is satisfied. Under conditions to be specified, the local search procedure is started from a point in the cluster.

In the next three subsections, we will describe three methods that fit in this framework, but differ in the rule by which they select the points that are added to the cluster and in the corresponding termination criterion. This difference is mainly due to the different ways in which the methods exploit the fact that the reduced sample is known to be a subset of the original sample of uniformly distributed points. The first two of these methods already appeared in an earlier article [Boender et al. 1982] and hence will be described only briefly. The analysis of their properties, however, is new.

3.1. Density clustering

Analogously to Törn [Törn 1976], we will let the clusters in this approach correspond to the reduced sample points in a subset T_i , $i=0,1,2,\ldots$, of S of stepwise increasing volume, where T_0 equals the seed point of the cluster and $T_{i+1} \supseteq T_i$, $i=1,2,\ldots$ A cluster is terminated if in a step no points are added to the cluster. However, we will adjust Törn's method in three ways; the choice of the seed points, the shape of the sets T_i and the increase in volume of these sets in each step.

We first turn to the choice of the seed points. As we will see later, it is advantageous to choose a local minimum as the seed point. Therefore, the local minima in X^* are first used as seed points. If all local minima known have been used as a seed point already, and there are still reduced sample points that have to be clustered, then a local search is started in the unclustered reduced sample point \overline{x} with smallest function value. If the resulting local minimum x^* was already known, then \overline{x} is assigned to the cluster that was initiated by x^* , and again a local search is started from the unclustered reduced sample point with smallest function value. If the resulting local minimum was not yet known, then it is chosen as the next seed point.

Let us now consider the shape of the sets T_i . Recall, that the cluster is initiated by a local minimum x* and that it should correspond to $L_{x^*}(y_k^{(\gamma kN)})$. This suggests to let T_i correspond to $L_{x^*}(y)$ for stepwise increasing values of y. The actual sets $L_{x^*}(y)$ may be hard to construct, but since f is twice continuously differentiable, we can approximate these sets by the level sets $\tilde{L}(y)$ around x* that are defined by the second order approximation \tilde{f} of f around x*:

$$\tilde{f}(x) = f(x^*) + \frac{1}{2}(x - x^*)^{T} H(x^*)(x - x^*).$$
(22)

Hence, in step i we let T_i be the set $\{x \in S | (x-x^*)^T H(x^*)(x-x^*) \leq r_i^2\}$, for some r_i to be determined below, with $r_{i+1} > r_i$, $i=1,2,\ldots$ (An approximation of $H(x^*)$ may be obtained, for example, as a byproduct of a quasi-Newton local search procedure.) Finally, we derive the rate at which r_i should increase with i so as to ensure proper termination of a cluster. The probability that the cluster is terminated in step i, equals the probability that the set $A_i = \{x \in S \mid x \in T_i, x \notin T_{i-1}\}$ does not contain any reduced sample points. To determine this probability for the case that there are still unclustered reduced sample points in $L_{x^*}(y_k^{(\gamma kN)})$, i.e. the probability of erroneous termination, we assume that the sets $L_{x^*}(y)$, with $f(x^*) \leq y \leq y_k^{(\gamma kN)}$ can be properly approximated by ellipsoids, so that $T_i \in L_{x^*}(y_k^{(\gamma kN)})$. Given this assumption, the probability of erroneous termination in step i, say α_k , equals the probability that none of the kN original sample points is located in A_i . Using (8) it follows that

$$\alpha_{k} = (1 - m(A_{i})/m(S))^{kN}.$$
(23)

Let us choose $m(A_i)$, and hence r_i , such that the probability α_k that the cluster is terminated incorrectly in step i, decreases with increasing k. For instance, if, for some $\sigma > 0$, $m(A_i) = (m(S)\sigma \log kN)/kN$, then

$$\alpha_{k} = \left(1 - \frac{\sigma \log kN}{kN}\right)^{kN}.$$
 (24)

It is not hard to verify that, for some constants c_1 , $c_2 > 0$, we have that

$$c_1 k^{-\sigma} \leq \left(1 - \frac{\sigma \log k}{k}\right)^k \leq c_2 k^{-\sigma}$$
(25)

for all k. Hence, if we terminate the cluster in step i if no unclustered reduced sample point exists in T_i with

$$r_{i} = \pi^{-\frac{1}{2}} \left(i\Gamma(1 + \frac{n}{2}) (\det H(x^{*}))^{\frac{1}{2}} m(S) \frac{\sigma \log kN}{kN} \right)^{1/n}, \quad (26)$$

then the probability that the cluster is terminated incorrectly in step i, decreases polynomially fast with increasing k.

A stepwise description of this method follows. Let we be the number of local minima x* with $f(x^*) \leq y_k^{(\gamma kN)}$, which are known at the start of the procedure.

Density Clustering

- Step 1. (Determine reduced sample) Determine the reduced sample by taking the γkN points with the smallest function values. Set j := 1.
- <u>Step 2</u>. (Determine seed points) Set i := 1. If all reduced sample points have been assigned to a cluster, stop. If $j \leq w$, then choose the j-th local minimum in X^* as the next seed point. If j > w, then apply P to the unclustered reduced sample point \overline{x} with the smallest function value. If the resulting local minimum x^* is an element of X^* , then assign \overline{x} to the cluster initiated by x^* and repeat step 2. If $x^* \in X^*$, then add x^* to X^* , set w := w+1 and let x^* be the next seed point.
- <u>Step 3.</u> (Form cluster) Add all unclustered reduced sample points which are within distance r_i of the seed point x* to the cluster initiated by x*. If no point has been added to the cluster for this specific value of r_i , then set j:= j+1 and go to Step 2, else set i := i+1 and repeat Step 3.

Unfortunately, if the set $L_{x^*}(y_k^{(\gamma kN)})$ differs substantially from an ellipsoid, then this influences both the probability that the cluster is terminated incorrectly and the probability that the cluster is expanded incorrectly, in an unpredictable way. To arrive at a satisfactory clustering method, it is necessary that the shape of the resulting clusters is not fixed. Intuitively the shape of the clusters should converge to the shape of the actual sets $L_{x^*}(y_k^{(\gamma kN)})$ with increasing k. A method which satisfies this property is presented in the next subsection.

3.2. Single Linkage Clustering

In the adapted Single Linkage method, the clusters are formed sequentially, and each cluster is again initiated by a seed point. After a cluster C is initiated, we find an unclustered point x such that

$$d(x,C) = \min_{\substack{x_1 \in C}} \|x - x_1\|$$
(27)

is minimal. We add x to C and repeat until d(x,C) exceeds the critical distance r_k .

Early implementations of Single Linkage and Density Clustering were the subject of limited computational experiments. These experiments showed that Single Linkage indeed approximates the sets $L_{x*}(y_k^{(\gamma kN)})$ more accurately than Density Clustering. However, to prove rigorously the superiority of Single Linkage it turns out that we must slightly adjust the rule according to which the seed points are selected. The reason is that it will turn out to be difficult to analyze Single Linkage in the regions near the boundary of S and in neighbourhoods of the elements of X*. Therefore, we will define the procedure so as never to start a local search in these regions. This may imply that no local search is started from any point in a certain cluster. This, however, is not a serious drawback if we first redefine S in a slightly different way. For some $\tau > 0$, we let \boldsymbol{Q}_{τ} be the set of points in S that are within distance τ of a point on the boundary of S, and we let \boldsymbol{S}_{τ} be the set of points in S which are not within distance T of a point on the boundary of S, so that $S_{\tau} = S \setminus Q_{\tau}$. We assume that all local minima of f occur in the interior of S_{τ} .

We will also have to give special treatment to the neighbourhoods of the elements of X^* . For some fixed and small υ , let X_{υ}^* be the set $\{x \in S \mid \|x - \overline{x}\| < \upsilon$, for any $\overline{x} \in X^*$. Recall that we assumed that a positive constant ε can be specified, such that the distance between any two stationary points exceeds ε . Hence, we can choose υ such that the distance between any two stationary points exceeds 2υ . We will now give a stepwise description of the adjusted Single Linkage procedure.

Single Linkage

- Step 1. (Determine reduced sample) Determine the reduced sample by taking the γkN sample points with the smallest function values. Let X^1 be the set of minima in X^* , let w be the number of elements of X^1 , and set j := 1.
- <u>Step 2</u>. (Determine seed points) If all reduced sample points have been assigned to a cluster, stop. If $j \leq w$, then choose the j-th local minimum in X^1 as the next seed point; go to Step 3. Determine the point \overline{x} which has the smallest function value among

the unclustered reduced sample points; \overline{x} is the next seed point. If $\overline{x} \in S_{\tau}$ and if $\overline{x} \in X_{\eta}^{*}$, then apply P to \overline{x} to find a local minimum x*; add new stationary points encountered during this search (possibly including x^*) to X^* .

Step 3 (Form cluster) Initiate a cluster by the seed point which is determined in Step 2. Add reduced sample points which are within distance r_{i} of a point already in the cluster to the cluster, until no more such points exist. Set j := j+1, and go to Step 2.

Let us now analyze Single Linkage, and determine an appropriate value for the critical distance rk. This critical distance will be chosen to depend on kN only so as to minimize the probabilities of two possible failures of the method: the probability that a local search is started, although the resulting minimum is known already, and the probability that no local search is started in a component of $L(y_k^{(\gamma kN)})$ which contains reduced sample points.

Let us first consider the probability that P is applied incorrectly to some reduced sample point. For a suitable choice of rk, we will prove that the probability that a local search is started, let alone started incorrectly, tends to 0 with increasing k. For this purpose we divide S into three subsets. Let Y_{ij} be the set of elements in S that are within distance υ of a stationary point of f. We already defined $Q^{}_{\tau}$ to be the set of elements in S that are within distance τ of a point on the boundary of S. Finally we let $M_{T,U}$ consist of the elements in S that do not belong to Q_{τ} or Y_{U} . Note that we defined Q_{τ} and Y_{U} as open sets, so that M_{τ} , is closed and therefore compact. We will start our analysis by considering the elements of $M_{\tau, \upsilon}$; the large majority of the reduced sample points belongs to this set. We wish to determine the probability that P is applied to a reduced sample point <u>x</u> with <u>x</u> = a $\in M_{\tau, \upsilon}$. Let $B_{a,r}$ be the set $\{x \in S \mid \|x-a\| \leq r\}$. Suppose that B contains a sample point z with f(z) < f(a). Clearly, z then belongs to the reduced sample, and if z is assigned to a cluster then a will be assigned to that cluster too. Moreover, it is easy to check that we will not apply the local search procedure to a before z has been assigned to a cluster. Thus, the probability that a local search is started in a reduced sample point $\underline{x} = a \in M_{\tau, \upsilon}$, is certainly smaller than the probability that there is no

sample point \underline{z} in B_{a,r_k} with $f(\underline{z}) < f(a)$. To calculate this latter probability, we need the following theorem.

THEOREM 7. For any $\tau > 0$ and $\upsilon > 0$, let a be an element of $M_{\tau,\upsilon}$, let $B_{a,r} = \{x \in S \mid \|x-a\| \leq r\}$, and let $A_{a,r} = \{x \in S \mid \|x-a\| \leq r \text{ and } f(x) \leq f(a)\}$. Then, uniformly in a,

$$\lim_{r \neq 0} \frac{m(A_{a,r})}{m(B_{a,r})} \geq \frac{1}{2}.$$
(28)

PROOF. Consider the set

$$D_{a,r} = \{x \in S \mid \|x-a\| \leq r \text{ and } g(a)^{T}(x-a) + \frac{1}{2}cr^{2} < 0\}, \quad (29)$$

where c is a positive constant which is greater than the supremum over S of the eigenvalues of H(x). From the Taylor expansion of f around a, we know that for all $x \in S$, with $||x-a|| \leq r$, there exists a θ , $0 \leq \theta \leq 1$, such that

$$f(x) - f(a) = g(a)^{T}(x-a) + \frac{1}{2}(x-a)^{T}H(a+\theta(x-a))(x-a)$$

$$\leq g(a)^{T}(x-a) + \frac{1}{2}c(x-a)^{T}(x-a)$$

$$\leq g(a)^{T}(x-a) + \frac{1}{2}cr^{2}.$$
(30)

Hence, if $x \in D_{a,r}$, then $x \in A_{a,r}$. Thus, we have proved that $D_{a,r} \subset A_{a,r}$.

Now consider an orthogonal matrix U (so that $U^{T}U = UU^{T} = I$), for which

$$U^{T}e_{1} = \frac{1}{\|g(a)\|} g(a),$$
 (31)

where e_1^T is the n-dimensional vector $(1,0,\ldots,0)$. Obviously, such a matrix U always exists, because condition (31) only fixes the first row of U to be equal to g(a)/||g(a)||, the norm of which is 1.

We can now rewrite the set D_{a,r} as follows

$$D_{a,r} = \{x \in S | (x-a)^T U^T U(x-a) \leq r^2 \text{ and } g(a)^T U^T U(x-a) + \frac{1}{2}r^2 c < 0\}$$

$$= \{x \in S | (U(x-a))^{T} U(x-a) \leq r^{2} \text{ and } (Ug(a))^{T} U(x-a) + \frac{1}{2}r^{2}c < 0\}$$
$$= \{x \in S | (U(x-a))^{T} U(x-a) \leq r^{2} \text{ and } \|g(a)\| = \frac{1}{2}U(x-a) + \frac{1}{2}r^{2}c < 0\}.$$

(32)

25

Hence, the matrix U defines a 1-1 correspondence between the elements of $D_{a,r}$ and the elements of

$$G_{a,r} = \{z \in \mathbb{R}^n | \|z\| \leq r \text{ and } \|g(a)\|e_1^T z + \frac{1}{2}r^2 c < 0\}.$$
 (33)

Note that for r sufficiently small, all points $x \in \mathbb{R}^n$ satisfying $\|x-a\| \leq r$ are contained in S, because it follows from the definition of $M_{\tau, \upsilon}$ that this is certainly true if $r < \tau$. The transformation defined by U does not change the distances between points and the angles between vectors, because, for every $x_1, x_2 \in \mathbb{R}_n$, $x_1^T x_2 = x_1^T U^T U x_2 = (Ux_1)^T U x_2 = z_1^T z_2$. Moreover, $m(D_{a,r}) = m(G_{a,r})$, since the determinant of U is 1. Since, for $r < \tau$, $B_{a,r} = \{x \in \mathbb{R}^n \mid \|x-a\| \leq r\}$, we obtain that $m(B_{a,r}) = r^n \pi^{n/2} / \Gamma(1 + \frac{n}{2})$. Since f is continuously differentiable, and $M_{\tau,\upsilon}$ is compact, the minimum of $\|g(x)\|$ over $M_{\tau,\upsilon}$ exists. Let p be this minimum; p cannot be zero because $M_{\tau,\upsilon}$ does not contain any stationary point, so that p > 0. Hence, if the first coordinate of z, say $z^{(1)}$, is smaller than $-\frac{1}{2}r^2 c/p$, and if $\|z\| \leq r$, then $z \in G_{a,r}$. Since the intersection of $\{z \in \mathbb{R}^n \mid \|z\| \leq r\}$ with the hyperplane $z^{(1)} = 0$ is a (n-1) dimensional hyperball with measure $\pi^{(n-1)\frac{1}{2}}r^{n-1}/\Gamma(1 + \frac{n-1}{2})$, it follows that

$$m(G_{a,r}) \ge \frac{1}{2} \cdot \frac{r^{n} \pi^{n/2}}{\Gamma(1+\frac{n}{2})} - \frac{r^{n-1} \pi^{\frac{1}{2}(n-1)}}{\Gamma(1+\frac{n-1}{2})} \cdot \frac{cr^{2}}{2p}$$
 (34)

Thus,

$$\lim_{\mathbf{r} \to 0} \frac{\mathbf{m}(\mathbf{A}_{a,\mathbf{r}})}{\mathbf{m}(\mathbf{B}_{a,\mathbf{r}})} \geq \lim_{\mathbf{r} \to 0} \frac{\mathbf{m}(\mathbf{D}_{a,\mathbf{r}})}{\mathbf{m}(\mathbf{B}_{a,\mathbf{r}})} = \lim_{\mathbf{r} \to 0} \frac{\mathbf{m}(\mathbf{G}_{a,\mathbf{r}})}{\mathbf{m}(\mathbf{B}_{a,\mathbf{r}})} \geq$$
(35)

$$\geq \lim_{r \neq 0} \frac{1}{2} - \frac{\operatorname{cr}^{n+1} \frac{1}{\pi^2} (n-1)}{2 p \Gamma(1 + \frac{n-1}{2})} \cdot \frac{\Gamma(1 + \frac{n}{2})}{r^n \pi^{n/2}} = \frac{1}{2}$$

Since the above reasoning is independent of the choice of a $\in M_{\tau,\upsilon}$ the result is now immediate.

Actually, the limit considered in (28) is precisely equal to $\frac{1}{2}$. This can be easily seen from the fact that we can still prove the Theorem (using a similar argument) if we change the definition of $A_{a,r}$ into $A_{a,r} = \{x \in S \mid \|x-a\| \leq r \text{ and } f(x) > f(a)\}.$

Theorem 7 is valid for any positive τ and positive υ , so that we can choose these numbers as small as we like. Note, that if $\upsilon = 0$, but a is not a stationary point, then the limit in (28) still equals $\frac{1}{2}$. However, the convergence is not uniform in a, because ||g(a)|| can become arbitrarily small.

Let us now return to the probability that, for some reduced sample point <u>x</u>, with <u>x</u> = $a \in M_{\tau, \upsilon}$, there exists a sample point <u>z</u> in B_{a, r_k} with $f(\underline{z}) < f(a)$, which bounds the probability that a local search is started in a. To calculate this probability, let us first consider the simpler case where <u>x</u> is an arbitrary sample point, with <u>x</u> = $a \in M_{\tau, \upsilon}$. The remaining kN-1 sample points are still distributed according to a uniform distribution over S and hence, the probability that none of these kN-1 uniform points is in A_{a,r_k} , i.e. is within distance r_k of a and has a smaller function value than a, equals (cf. (8))

$$(1 - m(A_{a,r_{k}})/m(S))^{kN-1}$$
. (36)

Moreover, provided that r_k tends to 0 with increasing k, we know from Theorem 7 that, for any β with $0 < \beta < \frac{1}{2}$, there exists a k_0 such that for $k > k_0$

$$\frac{m(A_{a,r_{k}})}{m(B_{a,r_{k}})} \geq \beta.$$
(37)

Hence, for any sample point $\underline{x} = a \in M_{\tau, \upsilon}$, the probability that there is no sample point \underline{z} in $B_{a, r_{\iota}}$ with $f(\underline{z}) < f(a)$ is smaller than

$$(1 - \beta m(B_{a,r_k})/m(S))^{kN-1}$$
 (38)

(for sufficiently large k).

Analogously to Subsection 3.1, we can choose r_k in such a way that probability (38) is constant or decreases with k. For instance, for some $\sigma > 0$, we can choose

$$r_{k} = \pi^{-\frac{1}{2}} \left(\Gamma(1 + \frac{n}{2}) m(S) \frac{\sigma \log k N}{k N} \right)^{1/n},$$
(39)

so that, for k large enough, $m(A_{a,r_{k}}) \geq (\beta \sigma \log kN)/kN$. Hence, for this specific choice of r_{k} we proved that the probability that for some sample point $\underline{x} = a \in M_{\tau,\upsilon}$, there is no sample point \underline{z} in $B_{a,r_{k}}$ with $f(\underline{z}) < f(a)$ is $O(k^{-\beta\sigma})$ (note that we can omit N in all O(.) terms since N is a constant). Since the number of sample points in iteration k is kN, we may conclude that the probability that there exists a sample point in $M_{\tau,\upsilon}$ which has no other sample point within distance r_{k} with smaller function value is $O(k^{1-\beta\sigma})$. It follows that the probability that there exists a reduced sample point in $M_{\tau,\upsilon}$ which has no other sample point in $M_{\tau,\upsilon}$ which has no other sample point in $M_{\tau,\upsilon}$ which has no ther sample point within distance r_{k} with smaller function value must also be $O(k^{1-\beta\sigma})$. Hence, we proved that, for any $\beta < \frac{1}{2}$, the probability that a local search is started from any element of $M_{\tau,\upsilon}$ in iteration k tends to 0 with increasing k.

Moreover, if we let $\underline{\xi}_k$ be the number of local searches started from points in $M_{\tau,\upsilon}$ in iteration k, and if we choose $\sigma > 4$, then it is easy to show that

$$\sum_{k=1}^{\infty} \Pr[\underline{\xi}_k > 0] < \infty.$$
(40)

Hence, it follows from the Borel-Cantelli Lemma that even if the sampling and clustering continues forever, then the total number of local searches ever started in $M_{T,U}$ is finite with probability 1.

We now turn to the probability that a local search is started in Q_{τ} and Y_{υ} . It follows from the description of Single Linkage that no local search will ever be started in Q_{τ} . To analyze the situation in Y_{υ} we need one more assumption: we assume that if we apply P to a point which is within distance υ of a stationary point \overline{x} , then we will recognize \overline{x} as such and add it to X^* (if necessary). Because we can choose υ as small as we

want, this assumption is reasonable. Hence, we start P at most once in the neighbourhood of any stationary point. Since the number of stationary points is finite, we may conclude that the probability that P is applied to a point in $Y_{\rm U}$ tends to 0 with increasing k

Thus, we proved the following theorem.

<u>THEOREM 8.</u> If the critical distance r_k of Single Linkage is determined by (39) with $\sigma > 2$, then the probability that a local search is started by Single Linkage in iteration k tends to 0 with increasing k. If $\sigma > 4$, then, even if the sampling continues forever, the total number of local searches ever started by Single Linkage is finite with probability 1.

We will now consider the second possible failure of Single Linkage, i.e. the possibility that no local minimum is found in a component of $L(y_k^{(\gamma kN)})$, although this component contains a sample point. We shall prove that with a probability increasing to 1, such a failure will not occur. In analyzing this probability, we again encounter the difficulty that the components of $L(y_k^{(\gamma kN)})$ depend on the specific value of the random variable $\underline{y}_k^{(\gamma kN)}$. As before, we therefore focus on the components of $L(y_{\gamma})$.

To examine the probability that no local minimum is found in a component of $L(y_{\gamma})$, say $L_{a}(y_{\gamma})$, although a sample point is contained in $L_{a}(y_{\gamma})$ we will first prove some general results. Roughly speaking, these results will show that if $x_{1} \in L_{a}(y_{\gamma})$ and $x_{2} \in L(y_{\gamma}) \setminus L_{a}(y_{\gamma})$, then the components of $L(y_{k}^{(\gamma kN)})$ containing x_{1} and x_{2} respectively are sufficiently far apart. The fact that those components have been defined to be closed sets will play an important role in the proofs. For any $y \in \mathbb{R}$ and any $a \in L(y)$, let $V_{a}(y)$ be the set of elements in L(y) which are not contained in $L_{a}(y)$. Furthermore, let the distance between two subsets E_{1} and E_{2} of the \mathbb{R}^{n} , say $d(E_{1},E_{2})$, be defined as the infimum of the distances between any element of E_{1} and any element of E_{2} .

<u>THEOREM 9.</u> For all $y \in \mathbb{R}$, there exists a $\delta > 0$ such that for all $a \in L(y)$, $d(L_a(y), V_a(y)) \geq \delta$.

<u>PROOF</u>. Since every component of L(y) contains a stationary point (see Theorem 3), L(y) only consists of a finite number of components. Hence, $V_a(y)$ is the union of a finite number of components and is therefore closed It follows from the definition of $L_a(y)$ and $V_a(y)$ that $L_a(y) \cap V_a(y) = \emptyset$. If $V_a(y) = \emptyset$, then the theorem is trivially true, so that we may assume that $V_a(y) \neq \emptyset$.

The remaining part of the proof is by contradiction; suppose that $d(L_a(y), V_a(y)) = 0$. Then there exist a sequence α_i in $L_a(y)$ and a sequence β_i in $V_a(y)$ with $\|\alpha_i - \beta_i\| < \frac{1}{i}$, $i = 1, 2, \dots$. Since both $L_a(y)$ and $V_a(y)$ are bounded (by S), both sequences contain a convergent subsequence, $\alpha_i(j) \stackrel{\text{and } \beta}{=} i(j)$, such that $\|\alpha_i(j) - \beta_i(j)\| \leq \frac{1}{i(j)}$, for every positive integer j, and

$$\lim_{j \to \infty} \alpha_{i(j)} = \alpha, \lim_{j \to \infty} \beta_{i(j)} = \beta.$$
(41)

Since both $L_a(y)$ and $V_a(y)$ are closed, we have that $\alpha \in L_a(y)$ and $\beta \in V_a(y)$ and $\|\alpha - \beta\| = 0$. This, however, contradicts $L_a(y) \cap V_a(y) = \beta$. Thus, there exist a δ_a such that $d(L_a(y), V_a(y)) > \delta_a$. Obviously, δ_a is equal for all a that belong to the same component of L(y). Since L(y)only consists of a finite number of components, we can choose δ independent of a, which completes the proof.

<u>THEOREM 10</u>. There exists an $\varepsilon > 0$ and a $\delta > 0$ such that for any $y \leq y_{\gamma} + \varepsilon$, for any $a \in L(y_{\gamma})$ and for any minimum $x^* \in L(y)$ which does not belong to $L_a(y_{\gamma})$, we have that $d(L_a(y), L_{x^*}(y)) \geq \delta$.

<u>PROOF</u>. If $y = y_{\gamma}$, then the result follows immediately from Theorem 9, since $L_{x*}(y) \subset V_a(y)$.

Now suppose that $y < y_{\gamma}$. It follows from the definition of a component that $L_a(y) \subset L_a(y_{\gamma})$ and that $L_{x^*}(y) \subset L_{x^*}(y_{\gamma})$. Hence, because of Theorem 9 there exists a $\delta_1 > 0$ such that

$$d(L_{a}(y),L_{x*}(y)) \geq d(L_{a}(y_{\gamma}),L_{x*}(y_{\gamma})) \geq d(L_{a}(y_{\gamma}),V_{a}(y_{\gamma})) \geq \delta_{1}.$$
(42)

The interesting case arises when $y > y_{\gamma}$. Since f only has a finite number of stationary points, there exists an ε_1 such that $L(y_{\gamma}+\varepsilon_1)$ contains no more stationary points than $L(y_{\gamma})$. Hence, we may assume that $x^* \in L(y_{\gamma})$. Suppose that a and x^* belong to the same component of $L(y_{\gamma}+\varepsilon)$ for any $\varepsilon > 0$. Then, they also belong to the same component of $A_i = \{x \in S \mid f(x) < y_{\gamma} + \frac{1}{i}\}$ for any positive integer i. It is easy to prove that A_i and its components are open. Hence, if a and x^* belong to the same component of A_i, then there exists a path joining a and x^* . Since a and x^* both belong to $L(y_{\gamma})$ and there exist a δ_1 , such that $d(L_a(y_{\gamma}), L_{x^*}(y_{\gamma})) \geq \delta_1$, we have that, for every i, there must exist an $\alpha_i \in A_i$ for which

$$y_{\gamma} < f(\alpha_{i}) < y_{\gamma} + \frac{1}{i},$$

$$d(\alpha_{i}, V_{a}(y_{\gamma})) > \frac{1}{3} \delta_{1},$$

$$d(\alpha_{i}, L_{a}(y_{\gamma})) > \frac{1}{3} \delta_{1}.$$
(43)

The sequence α_i contains a convergent subsequence $\alpha_{i(j)}$ such that $\lim_{j \to \infty} \alpha_{i(j)} = \alpha$ and $f(\alpha) = y_{\gamma}$. This, however, contradicts (43). We may conclude that there exists an ε such that a and x* do not belong to the same component of $L(y_{\gamma} + \varepsilon)$. Since there are only a finite number of components and a finite number of minima, we can choose ε independent of a and x*. By Theorem 9 it now follows that there exists a δ_2 (independent of a and x*) such that $d(L_a(y_{\gamma} + \varepsilon), L_{x*}(y_{\gamma} + \varepsilon)) \geq \delta_2$. Hence, if $y_{\gamma} < y \leq y_{\gamma} + \varepsilon$, then, for all $a \in L(y_{\gamma})$ and $x* \in L(y)$, $x* \notin L_a(y_{\gamma})$, we have proven that $d(L_a(y), L_{x*}(y)) \geq \delta_2$. By choosing $\delta = \min\{\delta_1, \delta_2\}$ the result is now immediate.

We now return to the probability that no local minimum is found by Single Linkage in a component $L_a(y_{\gamma})$ although there is a sample point in $L_a(y_{\gamma})$. We shall show first that $L_a(y_{\gamma})$ must contain a local minimum which is in a sense conveniently located. The possibility that $L_a(y_{\gamma})$ contains a number of local minima, only one of which may be discovered, creates some extra difficulties in the reasoning that follows below.

First, since $L_{a}(y_{\gamma})$ is compact, there must exist a point $e \in S$ which

is the global minimum of f over $L_a(y_{\gamma})$, i.e. f(e) < f(x) for all $x \in L_a(y_{\gamma})$. (We assume that the global minimum of f over $L_a(y_{\gamma})$ is unique; this, however, is not essential.) Since P cannot leave a component of a level set in which it is started (Theorem 3), and since e has the smallest function value in $L_a(y_{\gamma})$, it follows that if P is started in e then it stops in e as well. Hence, e is a local minimum of f over S, so that e is in the interior of S_{τ} .

If $f(e) = y_{\gamma}$, then the sample point in $L_a(y_{\gamma})$ must equal the local minimum e, so that a local minimum in $L_a(y_{\gamma})$ has been found.

To analyze the usual situation that $f(e) < y_{\gamma}$, it is convenient to prove the following theorem, which provides useful information about the location of the local minimum e.

<u>THEOREM 11</u>. If, for any component $L_a(y_{\gamma})$ of $L(y_{\gamma})$, e is the unique global minimum of f over $L_a(y_{\gamma})$, and if $f(e) < y_{\gamma}$, then there exists a neighbourhood E of e satisfying:

$$E \subset L_{a}(y_{\gamma}); \tag{44}$$

if $x_1 \in E$ and if $x_2 \in L_a(y_\gamma) \setminus E$, then $f(x_1) < f(x_2)$; (45)

if \overline{x} is any stationary point other than e, and if $x \in E$, then $||x - \overline{x}|| > v$; (46)

$$E \cap Q_{\tau} = \emptyset; \tag{47}$$

$$m(E) > 0 \tag{48}$$

<u>PROOF.</u> Let Y_{υ}^{e} be the set of points which are within distance υ of any stationary point other than e, and let $Z_{1} = \{x \in L_{a}(y_{\gamma}) \mid x \in Y_{\upsilon}^{e} \cup Q_{\tau}\}$. If $Z_{1} = \emptyset$, then we define \overline{y} to be y_{γ} , else \overline{y} is the infimum of f over Z_{1} . Now let E be the set $\{x \in L_{a}(y_{\gamma}) \mid f(x) < \overline{y}\}$. It is not hard to see that E satisfies (44)-(48).

It follows from (48) and (8) that the probability that E contains a sample point tends to 1 with increasing k. Suppose that $y_{\gamma}^{(\gamma kN)} \leq y_{\gamma} + \varepsilon$,

for the ε mentioned in Theorem 10 and suppose that E contains a reduced sample point x (the probability that both events occur simultaneously tends to 1 with increasing k). Let \overline{x} be the seed point of the cluster to which x is assigned. There are four possibilities:

 $\overline{x} \in E, \ \overline{x} \in L_{a}(y_{\gamma}) \setminus E, \ \overline{x} \in L(y_{\gamma}) \setminus L_{a}(y_{\gamma}) \text{ or } \overline{x} \notin L(y_{\gamma}).$

(i) If $\bar{x} \in E$, then it follows from (46) and (47) that either the local minimum e has been located already, or P is applied to \bar{x} to find a minimum in $L_a(y_{\gamma})$.

(ii) If $\overline{x} \in L_a(y_{\gamma}) \setminus E$, then $f(\overline{x}) > f(x)$ by (45). It follows from the description of Single Linkage that a point x cannot be assigned to a seed point \overline{x} which is not a minimum, if $f(\overline{x}) > f(x)$. Hence, \overline{x} must be a local minimum in $L_a(y_{\gamma})$.

(iii) Suppose that $\overline{x} \notin L(y_{\gamma})$. Since all seed points are in $L(y_{\gamma} + \varepsilon)$, and since there is no local minimum x^* with $y_{\gamma} < f(x^*) \leq y_{\gamma} + \varepsilon$ (see the proof of Theorem 10), it follows that \overline{x} is not a minimum. However, x cannot be assigned to a seed point \overline{x} which is not a minimum if $f(\overline{x}) > f(x)$. Hence \overline{x} cannot be outside $L(y_{\gamma})$, and this case cannot occur.

(iv) Suppose that $\bar{x} \in L(y_{\gamma}) \setminus L_{a}(y_{\gamma})$. Obviously, the component $L_{\bar{x}}(y_{\gamma})$ contains a minimum x^{*} , since if P is applied to \bar{x} , it converges to a minimum in $L_{\bar{x}}(y_{\gamma})$ by assumption. Hence, it follows from Theorem 10 that there is a $\delta > 0$, such that there is no point in $L_{a}(y_{\gamma})$ within distance δ of any point in $L_{\bar{x}}(y_{\gamma})$. It follows that if the critical distance r_{k} of Single Linkage is smaller than δ , then x cannot be assigned to the cluster initiated by \bar{x} .

Thus, if r_k tends to 0 with increasing k, and if $\underline{\xi}$ is the index of the iteration in which a local minimum is found in a component $L_a(y_{\gamma})$ which contains a sample point, then we have shown that the probability that $\underline{\xi}$ is less than k tends to 1 with increasing k. Hence, for every $\varepsilon > 0$ there exists a k_0 such that Pr $[\underline{\xi} < k_0] > 1-\varepsilon$, and we may conclude that $\underline{\xi}$ is finite with probability 1.

THEOREM 12. If the critical distance r_k of Single Linkage tends to 0 with increasing k, then, in every component of $L(y_{\gamma})$ in which a point has been sampled, a local minimum will be found by Single Linkage within a finite number of iterations with probability 1.

32

Note that we can omit the provision that the component $L_a(y_{\gamma})$ of $L(y_{\gamma})$ must contain a sample point, if we assume that the measure of $L_a(y_{\gamma})$ is positive. If no local minimum x* exists with $f(x*) = y_{\gamma}$, then this latter assumption is satisfied for every component of $L(y_{\gamma})$.

3.3. Mode Analysis

Density clustering and Single Linkage are based on very simple properties of the uniform distribution. In Density Clustering, a cluster is expanded if a certain region contains a reduced sample point and in Single Linkage a reduced sample point is assigned to a cluster if it is within the critical distance from a point which has already been assigned to the cluster. In principle it should be possible to design superior methods by using the information of more than two sample points simultaneously. The mode analysis approach to clustering [Wishart 1969] is an example of an approach where several points are used simultaneously to determine the regions in which there is a high density of points to be clustered. However, the method proposed in [Wishart 1969] is not suitable for our purpose since it ignores much of the information available, like the fact that the reduced sample points are a subset of the unform sample. The technique proposed in [Spircu 1979] (based on [Parzen 1962]) allows one to estimate the distribution from which the reduced sample points are drawn, ignoring however that this distribution changes over time through its dependence on $L(y_k^{(\gamma kN)})$. This difficulty can be overcome [Ruygrok 1982], but the resulting method is cumbersome and inferior to the much simpler one presented below.

We shall describe a method in which S is partitioned into small hypercubes or <u>cells</u>. We say that a cell A is <u>full</u> if it contains more than

$$\frac{1}{2} \frac{\mathrm{m}(\mathrm{A})\mathrm{kN}}{\mathrm{m}(\mathrm{S})}$$
(49)

reduced sample points (i.e. more than half the expected number of sample points in A). If a cell is not full it is <u>empty</u>. We say that two cells are <u>neighbours</u>, or <u>neighbouring</u> cells, if they contain elements which are arbitrarily close to each other. We shall let a cluster correspond to a connected subset of S which corresponds to a number of full cells. These clusters can be found by applying a Single Linkage type algorithm to the full cells, such that if two cells are neighbours, then they are assigned to the same cluster.

A stepwise description follows below. (The sets Q_{τ}, Y_{υ} and X_{υ}^{*} are needed and defined for the same reasons as in the previous subsection.)

Mode Analysis

Step 1. (Determine reduced sample) Determine the reduced sample by taking the γkN sample points with the smallest function values.

Step 2. (Define cells) Divide S into v cells.

- Step 3. (Determine full cells) For each cell, determine the number of reduced sample points in the cell. If this number exceeds (49) then the cell is full, else it is empty.
- Step 4. (Determine seed cell) If all full cells have been assigned to a cluster, stop.

If an unclustered full cell exists which contains a minimum which is in X^* , then this cell is the new seed cell; go to Step 5. Determine the point \overline{x} which has the smallest function value among the reduced sample points which are in unclustered full cells. The cell which contains \overline{x} is the new seed cell. If $\overline{x} \in S$ and if $\overline{x} + x^*$ the new seed cell.

If $\overline{x} \in S_{\tau}$ and if $\overline{x} \notin X_{\upsilon}^{*}$, then apply P to \overline{x} to find a local minimum x*; add new stationary points encountered during this search (possibly including x*) to X^{*}.

Step 5. (Form cluster) A cluster is initiated by the seed cell which is

determined in Step 4.

Full cells which are a neighbour of a cell already in the cluster are assigned to the cluster, until there are no more such cells. Go to Step 4.

Since the properties of Mode Analysis do not really depend on $\sqrt[n]{v}$ being integer, we will, for the sake of analysis, assume that S is a hypercube and that $\sqrt[n]{v}$ is an integer so that S can be divided in v equal hypercubes. For some $\sigma > 0$, we choose v to be equal to $kN/(\sigma \log kN)$ so that each cell has measure $(m(S)\sigma \log kN)/kN$.

Intuitively speaking, we can say that Mode Analysis and Single Linkage will result in similar clusters and similar sets X^* , since the measure of the points within the critical distance (39) of a given point equals the measure of a cell. However, Mode Analysis seems somewhat less dependent on the particular irregularities of the sample, because it considers a number of sample points at the same time (σ logkN in expectation).

It is not possible to prove the superiority of either of the two methods rigorously. For instance, consider a one dimensional function, such that $L(y_k^{(\gamma kN)})$ consists of two components. Let the distance between both components be d and let r_k equal (m(S) σ logkN)/kN, i.e. the critical distance of Single Linkage and the cell size of Mode Analysis. Let us assume that the probability that a cell is full is high if the fraction of the cell that intersects with $L(y_k^{(\gamma kN)})$ exceeds $\frac{1}{2}$, and that a cell is probably empty if this fraction is smaller than $\frac{1}{2}$. (We shall see later that this is true if σ is large enough.) Obviously, if $d > 2r_k$, then both methods will recognize both components as such, and if $d < \frac{1}{2}r_k$, then it is likely that both methods will fail to detect both components. If $r_k < d < 2r_k$, then Single Linkage always detects both components, whereas, with small probability, Mode Analysis may fail. However, if $\frac{1}{2}r_k < d < r_k$, then the probability that Single Linkage will assign all reduced sample points to the same cluster is considerable, since this will happen if two reduced sample points exist in the different components which are within distance rk of each other. For Mode Analysis this probability is smaller, since it is possible that the region between both components covers an important part of one of the cells, in which case this cell is probably empty.

Clearly, the above arguments loose their relevance if k tends to infinity, since the distance between the components does not tend to zero with increasing k and r_k does. For k large enough it turns out that Mode Analysis and Single Linkage are very similar. Actually, it is possible to adjust the analysis of Single Linkage such that it can be applied to Mode Analysis to yield similar results. The most important adjustment of the analysis is needed in Theorem 7, for which we will now given the appropriate extension.

THEOREM 13. Let S be partitioned into equal hypercubes with edgelength r_1 . For any $\tau > 0$ and $\upsilon > 0$, let a be an element of $M_{\tau,\upsilon}$, and let C_1 be the cell which contains a. Then there exists a cell C_2 (depending on r_1) which is a neighbour of C_1 such that, uniformly in a

$$\lim_{r_1 \neq 0} \frac{m(\{x \in C_2 \mid f(x) < f(a)\})}{m(C_2)} = 1.$$
 (50)

<u>PROOF</u>. We will proof this theorem by adjusting the proof of Theorem 7; notations used in this latter proof have the same meaning here. We can choose r and r_1 , such that $r = 2r_1 \sqrt{n}$. Hence, we may assume that C_1 and all its neighbours are contained in $B_{a,r}$. We know that there exists a transformation z = U(x-a) which maps $B_{a,r}$ into $\{z \in \mathbb{R}^n \mid \|z\| \leq r\}$, such that the orthogonal matrix U satisfies (31) and such that the image of the hypercube C_1 is a hypercube with edgelength r_1 containing the origin.

We first prove that there exists a cell C_2 which is a neighbour of C_1 of which the image (under the above transformation) is completely contained in $A^- = \{z \in \mathbb{R}^n | z^{(1)} \leq 0\}$. For this purpose, note that each hypercube has a vertex of which the first coordinate is smaller than or equal to the first coordinate of any other member of the hypercube. Hence, C_1 has a vertex, say a_1 , which has the property that, for all $x \in C_1$,

$$e_1^{\mathrm{T}} \mathrm{U}(a_1 - a) \leq e_1^{\mathrm{T}} \mathrm{U}(x - a).$$
⁽⁵¹⁾

Since $a \in C_1$ and $e_1^T U(a-a) = 0$, it follows that the image of a_1 , $U(a_1-a)$ is in A⁻. Obviously, each vertex (in $M_{\tau, \upsilon}$) is shared by 2^n cells, and each of these cells can be characterized by the fact whether or not the i-th coordinate of the elements in the cell is greater than the i-th coordinate of this vertex $(i=1,2,\ldots,n)$. More formally, a cell which has a_1 as a vertex consists of elements x that can be written as

$$x = a_{1} + \sum_{i=1}^{n} \lambda_{i} e_{i}, \qquad (52)$$

where e_i is the i-th unit vector, and either $-r_1 \leq \lambda_i \leq 0$ or $0 \leq \lambda_i \leq r_1$ for every i=1,2,...,n. Now consider a cell, say C_2 , whose elements satisfy (52) where

$$\lambda_{i} \leq 0 \text{ if } e_{1}^{T} U e_{i} \geq 0,$$

$$\lambda_{i} \geq 0 \text{ if } e_{1}^{T} U e_{i} < 0.$$
(53)

Hence,

$$e_{1}^{T}U(x-a) = e_{1}^{T}U((a_{1}-a) + \sum_{i=1}^{n} \lambda_{i}e_{i})$$

$$= e_{1}^{T}U(a_{1}-a) + \sum_{i=1}^{n} \lambda_{i}e_{1}^{T}Ue_{i}$$

$$\leq 0.$$
(54)

It follows that C_2 is completely contained in A⁻. It is easy to show that C_2 is not C_1 since, for all $x \in C_2$ we have that

$$e_1^T U(x-a) \leq e_1^T U(a_1-a),$$
 (55)

while for the elements of C_1 the reverse is true. Thus, we have proved that C_2 is a neighbour of C_1 which is completely contained in A⁻.

We know from the proof of Theorem 7 that there exist positive constants c and p, such that the elements whose images are in

$$\{z \in \mathbb{R}^n \mid \|z\| \leq r \text{ and } z^{(1)} \leq -\frac{1}{2} \frac{r^2 c}{p}\}$$
 (56)

have a function value smaller than f(a). Since the image of C_2 is a hypercube with edgelength r_1 , which is completely contained in $\{z \in \mathbb{R}^n \mid \|z\| \leq r \text{ and } z^{(1)} \leq 0\}$, simple calculations yield

$$m(\{x \in C_2 | f(x) < f(a)\}) \ge r_1^n - (r_1 \sqrt{n})^{n-1} \cdot \frac{r^2 c}{2p}$$
 (57)

Thus,

$$\frac{\lim_{r_{1}^{+} 0} \frac{m(\{x \in C_{2} | f(x) < f(a)\})}{m(C_{2})}}{\lim_{r_{1}^{+} 0} \frac{r_{1}^{n} - r_{1}^{n-1} \cdot \frac{r^{2} c n^{\frac{1}{2}(n-1)}}{2p}}{r_{1}^{n} 0 - \frac{r_{1}^{n} r_{1}^{n} \cdot \frac{r_{2}^{n}}{p}}{r_{1}^{\frac{1}{2}(n+1)}} = \lim_{r_{1}^{+} 0} 1 - \frac{2r_{1} c n^{\frac{1}{2}(n+1)}}{p} = 1.$$
(58)

(Recall that $r = 2r_1 / n_{\cdot}$) Since the above arguments are independent of a, the result is now immediate.

To determine the probability that a local search is started by Mode Analysis, we divide S in the three sets Q_{τ} , Y_{υ} and $M_{\tau,\upsilon}$ again. As in Single Linkage, no local search is ever started in Q_{τ} , and, from our earlier assumption, for any stationary point, only one local search can be started within distance υ of this point.

Now let us consider the probability that P is applied to a reduced sample point $\underline{x} = a \in M_{\tau, \psi}$. Obviously, Mode Analysis will not start P from a point a if the cell containing a has a neighbour which is full and contains a sample point with a function value smaller than f(a). To determine the probability of the above event, let us first consider the simpler case where a is an arbitrary sample point in $M_{\tau, \psi}$.

Let C_1 be the cell in which a is located. Obviously, the remaining kN-1 sample points are still distributed according to a uniform distribution over S. From Theorem 13 we know, that for any $\frac{1}{2} < \beta < 1$, there exists a k_0 such that if $k > k_0$, then there exists a neighbour C_2 of C_1 with

$$\frac{m(\{x \in C_2 | f(x) < f(a)\})}{m(C_2)} \ge \beta.$$
(59)

Since $m(C_2) = (m(S)\sigma \log kN)/kN$, the probability that less than $\frac{1}{2}\sigma \log kN$ of the kN-1 sample points are in $\{x \in C_2 | f(x) < f(a)\}$ is smaller than

$$\sum_{i=0}^{\frac{1}{2}\sigma \log kN-1} {\binom{kN-1}{i}} \left(\frac{\beta\sigma \log kN}{kN}\right)^{i} \left(1 - \frac{\beta\sigma \log kN}{kN}\right)^{kN-i-1}.$$
 (60)

Obviously, for k large enough (60) is smaller than

$$\sum_{i=1}^{\frac{1}{2}\sigma \log kN} 2\binom{kN}{i} \left(\frac{\beta\sigma \log kN}{kN}\right)^{i} \left(1 - \frac{\beta\sigma \log kN}{kN}\right)^{kN-i}.$$
 (61)

Using Chernoff's inequality [Erdös & Spencer 1974], it follows that we can chose β such that (61) is $0(k^{-\sigma/10})$. Hence, for an arbitrary sample point a in $M_{\tau, \upsilon}$, we proved that the probability that the cell containing a has no neighbour with more than $\frac{1}{2}\sigma \log N$ sample points with a function value smaller than f(a) is $0(k^{-\sigma/10})$. (Of course, this is not the sharpest possible bound, but it suffices for our purpose.) Since the number of sample points in iteration k is kN, the probability that there exists a sample point a in $M_{\tau, \upsilon}$, such that the cell containing a has no neighbour with more than $\frac{1}{2}\sigma \log N$ sample points with a function value smaller than f(a) is $0(k^{1-\sigma/10})$. It is not difficult to verify that the same statements are true with respect to the reduced sample. Hence, if $\sigma > 10$, then the probability that a local search is started by Mode Analysis in iteration k tends to 0 with increasing k. As in Single Linkage, we can use the Borel-Cantelli Lemma to prove the following analagon of Theorem 8.

THEOREM 14. If the number of cells in Mode Analysis is kN/($\sigma \log N$) with $\sigma > 10$, then the probability that a local search is started by Mode Analysis in iteration k tends to 0 with increasing k. If $\sigma > 20$, then, even if the sampling continues forever, the total number of local searches started by Mode Analysis is finite with probability 1.

Let us now consider the second possible failure of Mode Analysis, i.e. the possibility that no local minimum is found in a component $L_a(y_{\gamma})$ although a sample point exists in $L_a(y_{\gamma})$. The analysis of this possibility is similar to the analysis of the corresponding possibility for Single Linkage. Ignoring details, we just state the final result.

39

<u>THEOREM 15.</u> If the number of cells in Mode Analysis is kN/(σ logkN) with $\sigma > 0$, then, in every component L(y_{γ}) in which a point has been sampled, a local minimum will be found by Mode Analysis within a finite number of iterations with probability 1.

(62)

4. CONCLUDING REMARKS.

The methods that have been described in the three previous subsections share one major deficiency. Although we know that a region of attraction cannot intersect with two different components of a level set, it is possible that a component of $L(y_k^{(\gamma kN)})$ contains more than one region of attraction. Since only one local search is started in each cluster, it is therefore possible that a local minimum may not be found although its region of attraction contains a reduced sample point. In this section we briefly consider three possible remedies to overcome this problem.

A first remedy is to replace every reduced sample point x by another point which is the result of a <u>steepest descent step</u> started in x, i.e. a one dimensional search from x in the direction of the negative gradient in x. The clustering procedure can then be applied to the resulting <u>transformed sample</u> as though it was the reduced sample [Boender et al. 1982]. From a theoretical point of view, however, the transformed sample has the disadvantage that its elements are no longer a subset of the original uniform sample. Thus, the analysis of the clustering methods that has been described in the previous subsections is no longer valid.

A second remedy is based on the observation that the negative gradient at a point which belongs to the region of attraction of a minimum x^* , will generally have a component in the direction of x^* . The methods described in the foregoing subsections can be improved by inspecting the gradient at a point before assigning it to a cluster. More precisely, if a cluster is initiated by a seed point \overline{x} (usually a local minimum), we then approximate the derivative of f in x in the direction of \overline{x} by

 $\frac{f(x+h(x-x))-f(x)}{h\|x-x\|}$

for small h, and we reject x for the cluster if this value is positive. In the case of Mode Analysis, we could inspect the gradient at the reduced sample point with the smallest function value in a cell before assigning the cell to the cluster. Although this <u>gradient criterium</u> can be incorrect and may affect the methods in an unpredictable way, it turns out to be very useful from a computational point of view. (see [Boender et al. 1982]).

A third possible remedy affects the methods even more deeply. Since we are only interested in the global minimum, we could reduce the sample even further. If, for some $\gamma \in (0,1)$, we would (re)define the reduced sample to contain the $(kN)^{\gamma}$ sample points with the smallest function values, we could prove much stronger (asymptotic) results. However, the statements are valid only if all reduced sample points are arbitrary close to a global minimum. Obviously, at that moment the problem has been solved a long time ago already. Moreover, the analysis would not yield any insight into the way in which the resulting methods would function before we arrive in the asymptotic case.

We conclude that the idea of sample reduction and clustering gives rise to interesting methods, but does not solve the problem satisfactorily from a theoretical point of view. In particular we cannot always avoid that a cluster contains several regions of attraction, so that a (local) minimum may still be missed. In Part II of this paper we will deal with this problem in a more fundamental way.

ACKNOWLEDGEMENTS

The research of the first author was partially supported by a NATO Scientist Fellowship. The research of the second author was partially supported by the Netherlands Foundation for Mathematics SMC with financial aid from the Netherlands Organization for Advancement of Pure Research (ZWO).

REFERENCES

Amsterdam).

Anderssen, R.S. (1972), Global optimization. In R.S. Anderssen, L.S. Jennings and D.M. Ryan (eds.), <u>Optimization</u> (University of Queensland Press).

Anderssen, R.S. and P. Bloomfield (1975), Properties of the random search in global optimization. Journal of Optimization Theory and Applications 16, 383-398.

Archetti, F. and B. Betro (1978a), A priori analysis of deterministic strategies for global optimization. In [Dixon & Szegö 1978a].

Archetti, F. and B. Betro (1978b), On the effectiveness of uniform random sampling in global optimization problems. Technical Report, University of Pisa, Pisa, Italy.

Archetti, F. and F. Frontini (1978), The application of a global optimization method to some technological problems. In [Dixon & Szegö 1978a].

Bahadur, R.R. (1966), A note on quantiles in large samples. <u>Annals of</u> <u>Mathematical Statistics 37</u>, 577-580.

Becker, R.W. and G.V. Lago (1970), A global optimization algorithm. In <u>Proceedings of the 8th Allerton Conference on Circuits and Systems</u> <u>Theory</u>.

Betro, B. (1981), Bayesian testing of nonparametric hypotheses and its application to global optimization. Technical Report, CNR-IAMI, Italy.

Boender, C.G.E., A.H.G. Rinnooy Kan, L. Stougie and G.T. Timmer (1980), Global optimization: a stochastic approach. In F. Archetti and M. Cugiani (eds.), <u>Numerical Techniques for Stochastic Systems</u> (North-Holland, Amsterdam).

Boender, C.G.E., A.H.G. Rinnooy Kan, L. Stougie and G.T. Timmer (1982), A stochastic method for global optimization. <u>Mathematical Programming</u> <u>22</u>, 125-140.

Boender, C.G.E. & A.H.G. Rinnooy Kan (1983), 'A Bayesian analysis of the number of cells of a multinomial distribution', <u>The Statistician 32</u>.

Boender, C.G.E. (1984), <u>The Generalized Multinomial Distribution:</u> <u>A Bayesian Analysis and Applications</u>. Ph.D. Dissertation, Erasmus Universiteit Rotterdam (Centrum voor Wiskunde en Informatica,

Boender, C.G.E. & A.H.G. Rinnooy Kan (1985), 'Bayesian Stopping rules for a class of stochastic global optimization methods', Technical Report, Econometric Institute, Erasmus University Rotterdam.

- Boender, C.G.E., A.H.G. Rinnooy Kan and G.T. Timmer (1985), 'A stochastic approach to global optimization'. To appear in <u>Proceedings of the NATO</u> <u>ASI on Computational Mathematical Programming</u>, Bad Windsheim.
- Brooks, S.H. (1958), A discussion of random methods for seeking maxima. Operations Research 6, 244-251.
- Chung, K.L. (1974), <u>A course in Probability Theory</u> (Academic Press, London).
- Devroye, L. (1978), Progressive global random search of continuous functions. <u>Mathematical Programming</u> 15, 330-342.
- Dixon, L.C.W. and G.P. Szegö (eds.) (1975), <u>Towards Global Optimization</u> (North-Holland, Amsterdam).
- Dixon, L.C.W., J. Gomulka and G.P. Szegö (1975), Towards global optimization. In [Dixon & Szegö 1975].
- Dixon, L.C.W. and G.P. Szegö (eds.) (1978a), <u>Towards Global Optimization 2</u> (North-Holland, Amsterdam).
- Dixon, L.C.W. and G.P. Szegö (1978b), The global optimization problem. In [Dixon & Szegö 1978a].
- Erdös, P. and J. Spencer (1974), <u>Probabilistic Methods in Combinatorics</u> (Academic Press, London).
- Hartman, J.K. (1973), Some experiments in global optimization. <u>Naval</u> <u>Research Logistics Quarterly 20, 569-576.</u>
- Ivanov, V.V. (1972), On optimal algorithms of minimization in the class of functions with the Lipschitz condition. Information Processing 2, 1324-1327.
- Parzen, E. (1962), On estimation of a probability density functions and mode. <u>Annals of Mathematical Statistics 33</u>, 1065-1076.
- Rinnooy Kan, A.H.G. and G.T. Timmer (1984), Stochastic methods for global optimization. To appear in the <u>American Journal of Mathematical and</u> <u>Management Sciences</u>.
- Rinnooy Kan, A.H.G. & G.T. Timmer (1985), 'Stochastic global optimization methods. Part II: multi level methods'.
- Rubinstein, R.Y. (1981), <u>Simulation and the Monte Carlo Method</u> (John Wiley & Sons, New York).

Ruygrok, A.J. (1982), <u>Mode Analysis in Globaal Optimaliseren</u>. Master Thesis, Erasmus University, Rotterdam (in Dutch).

Sobol, I.M. (1982), On an estimate of the accuracy of a simple multidimensional search. <u>Soviet Math. Dokl. 26</u>, 398-401.

- Solis, F.J. and R.J.E. Wets (1981), Minimization by random search techniques. <u>Mathematics of Operations Research 6</u>, 19-30.
- Spircu, L. (1979), Cluster analysis in global optimization. <u>Economic</u> <u>Computation and Economic Cybernetic Studies and Research 13, 43-50.</u>

Sukharev, A.G. (1971), Optimal strategies of the search for an extremum. Computational Mathematics and Mathematical Physics 11, 119-137.

Törn, A.A. (1976), Cluster analysis using seed points and density determined hyperspheres with an application to global optimization. In <u>Proceeding of the third International Conference on Pattern</u> <u>Recognition</u>, Coronado, California.

Törn, A.A. (1978), A search clustering approach to global optimization. In [Dixon & Szegö 1978a].

- Wishart, D. (1969), Mode Analysis: a generalization of nearest neighbour which reduces chaining effects. In A.J. Cole (ed.), <u>Numerical Taxonomy</u> (Academic Press, New York).
- Wolfe, P. (1969), Convergence conditions for ascent methods. <u>Siam Review</u> <u>11</u>, 226-235.

Wolfe, P. (1971), Convergence conditions for ascent methods II: some corrections. <u>Siam Review 13</u>, 185-188.

Zielinski, R. (1981), A stochastic estimate of the structure of multiextremal problems. <u>Mathematical Programming 21</u>, 348-356.

LIST OF REPORTS 1985

- 8500 "Publications of the Econometric Institute Second Half 1984: List of Reprints 378-400, Abstracts of Reports".
- 8501/0 R.M. Karp, J.K. Lenstra, C.J.H. McDiarmid and A.H.G. Rinnooy Kan, "Probabilistic analysis of combinatorial algorithms: an annotated bibliography", 26 pages.
- 8502/E M.E. Homan, "Verschillen in consumptie tussen één en tweekostwinnerhuishoudens: een eerste analyse", 13 pages.
- 8503/0 A.W.J. Kolen, "The round-trip p-center and covering problem on a tree", 17 pages.
- 8504/0 H.C.P. Berbee, C.G.E. Boender, A.H.G. Rinnooy Kan, C.L. Scheffer,
 R.L. Smith and J. Telgen, "Hit-and-run algorithms for the identification of nonredundant linear inequalities", 32 pages.
- 8505/S H. Brozius and L. de Haan, "On limiting laws for the convex hull of a sample", 10 pages.
- 8506/0 J.B.G. Frenk and A.H.G. Rinnooy Kan, "On the rate of convergence, to optimality of the LPT rule postscript", 5 pages.
- 8507/E P. Kooiman, H.K. van Dijk and A.R. Thurik, "Likelihood diagnostics and Bayesian analysis of a micro-economic disequilibrium model for retail services", 35 pages.
- 8508/0 C.G.E. Boender, A.H.G. Rinnooy Kan and J.R. de Wit, "A Bayesian procedure for the (s,Q) inventory problem", 26 pages.
- 8509/E B.M.S. van Praag, S. Dubnoff and N.L. van der Sar, "From judgments to norms: measuring the social meaning of income, age and educaton", 35 pages.
- 8510/S B. Bode, J. Koerts and A.R. Thurik, "On shopkeepers' pricing behaviour", 19 pages.
- 8511/M H. Bart, I. Gohberg and M.A. Kaashoek, "Exponentially dichotomous operators and inverse Fourier transforms", 70 pages.
- 8512/M J. Brinkhuis, "Normal integral bases and embedding of fields", 13 pages.
- 8513/E P.M.C. de Boer, "On the relationship between Revankar's and Lu-Fletcher's production function", 5 pages.
- 8514/S L. de Haan, "Extremes in higher dimensions: the model and some statistics", 15 pages.
- 8515 Publications of the Econometric Intitute First Half 1985: List of Reprints 401-414, Abstracts of Reports.

- 8516/A A.R. Thurik and A. Kleijweg, "Cyclical effects in retail labour productivity", 20 pages.
- 8517/C B.M.S. van Praag and J. van Weeren, "The impact of past experiences and anticipated future on individual income judgements", 24 pages.
- 8518/A A.R. Thurik and J. Koerts, "Behaviour of retail entrepreneurs". 16 pages.
- 8519/B M. Hazewinkel, J.F. Kaashoek and B. Leynse, "Pattern formation for a one dimensional evolution equation based on Thom's River Basin Model". 24 pages.
- 8520/A H.K. van Dijk, T. Kloek and C.G.E. Boender, "Posterior moments computed by mixed integration", 25 pages.
- 8521/B J. Brinkhuis, "Testing concavity and quasi-concavity is easy", 12 pages.
- 8522/A R. Harkema, "Minimum sample size requirements for maximum likelihood estimation of some demand models", 25 pages.
- 8523/B A.C.F. Vorst, "The general linear group of discrete Hodge algebras", 10 pages.
- 8524/A A.M.H. Gerards and A.W.J. Kolen, "Polyhedral combinatorics in combinatorial optimization", 25 pages.
- 8525/A B.J. Lageweg, J.K. Lenstra, A.H.G. Rinnooy Kan and L. Stougie, "Stochastic integer programming by dynamic programming", 19 pages.
- 8526/A J.B. Orlin, "A dual version of Tardo's algorithm for linear programming", 9 pages.
- 8527/A B.M.S. van Praag, "Household cost functions and equivalence scales", 24 pages (report 8424/E revised)
- 8528/A S. Schim van der Loeff, "Limited information maximum likelihood estimation of a subsystem of nonlinear equations", .. pages.
- 8529/A B.S. van der Laan and A.S. Louter, "On the number and the amount of damage of a passenger car traffic accidents in the Netherlands", ..pages.
- 8530/A B.S. van der Laan and A.S Louter, "A statistical model for the costs of passenger car traffic accidents", .. pages.
- 8531/A B.M.S. van Praag, J. de Leeuw and T. Kloek, "The population-sample decomposition approach to multivariate estimation methods", 34 pages.
- 8532/A M.E. Homan, B.M.S. van Praag and A.J.M. Hagenaars, "Household cost functions and the value of home production in one-and two earner families", .. pages.
- 8533/B R.J. Stroeker, "An inequality for YFF's analogue of the Brocard angel

- 8534/A J. Bouman, "Testing nonnested linear hypothesis: A Bayesian approach based on imcomfetely specified prior distributions". pages.
- 8535/A P.M.C. de Boer, R. Harkema and B.J. van Heeswijk, "Estimating foreign trade functions, a comment and a correction". pages.
- 8536/A C.G.E. Boender and A.H.G. Rinnooy Kan, "Bayesian stopping rules for multistart global optimization methods", 30 pages.
- 8537/A A.H.G. Rinnooy Kan and G.T. Timmer, "The multi level single linkage method for unconstrained and constrained global optimization", 14 pages.
- 8538/A A.H.G. Rinnooy Kan, "Probabilistic analysis of algorithms", 26 pages.
- 8539/A A.H.G. Rinnooy Kan and G.T. Timmer, "Stochastic global optimization methods, Part I: Clustering methods", 44 pages.

Until report 8515 this series was devided in E(conometrics), S(tatistics), M(athematics) and O(perations Research). From report 8516 on it will be

A. Economics, Econometrics and Operations Research

B. Mathematics

C. Miscellaneous

